

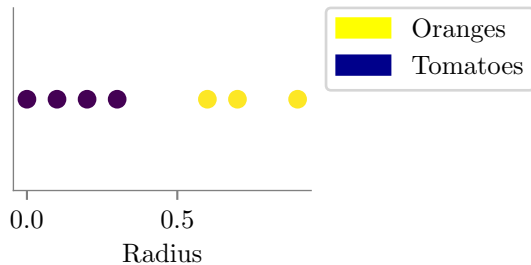
Logistic Regression

Nipun Batra

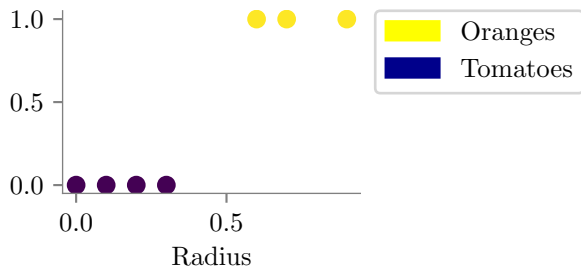
IIT Gandhinagar

August 1, 2025

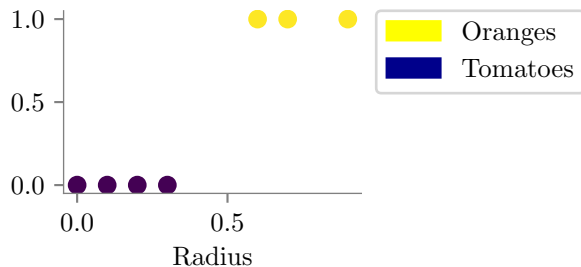
Classification Technique



Classification Technique

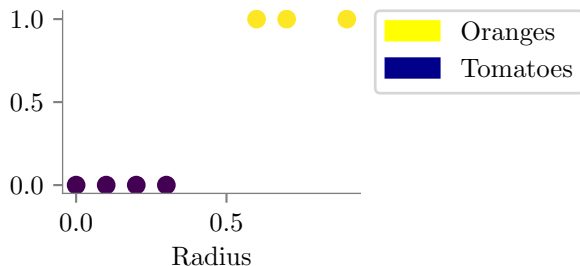


Classification Technique



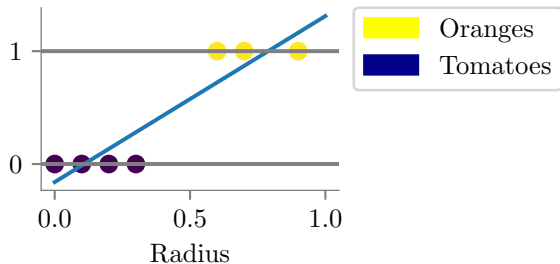
Aim: $\text{Probability}(\text{Tomatoes} \mid \text{Radius})$? or

Classification Technique



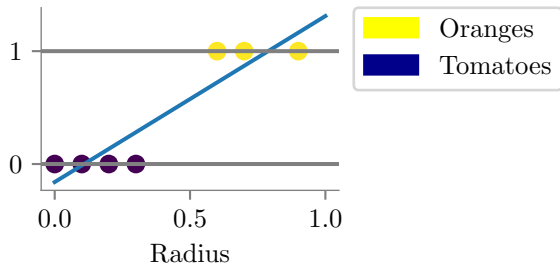
Aim: Probability(Tomatoes | Radius) ? or
More generally, $P(y = 1|X = x)$?

Idea: Use Linear Regression



$$P(X = \text{Orange} | \text{Radius}) = \theta_0 + \theta_1 \times \text{Radius}$$

Idea: Use Linear Regression



$$P(X = \text{Orange} | \text{Radius}) = \theta_0 + \theta_1 \times \text{Radius}$$

Generally,

$$P(y = 1 | x) = X\theta$$

Idea: Use Linear Regression

Prediction:

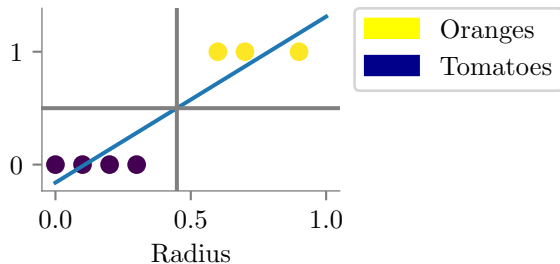
If $\theta_0 + \theta_1 \times \text{Radius} > 0.5 \rightarrow \text{Orange}$
Else $\rightarrow \text{Tomato}$

Problem:

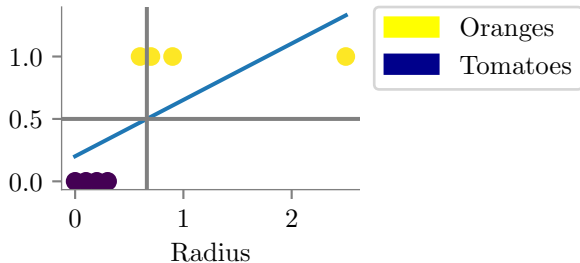
Range of $X\theta$ is $(-\infty, \infty)$

But $P(y = 1 | \dots) \in [0, 1]$

Idea: Use Linear Regression

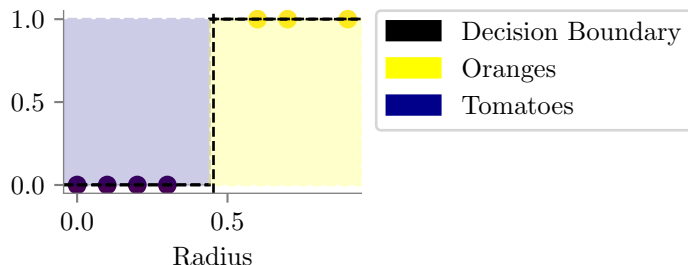


Idea: Use Linear Regression



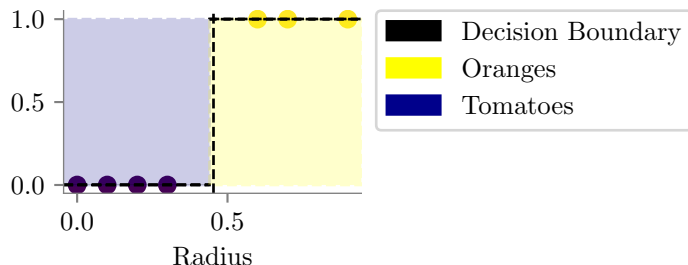
Linear regression for classification gives a poor prediction!

Ideal boundary



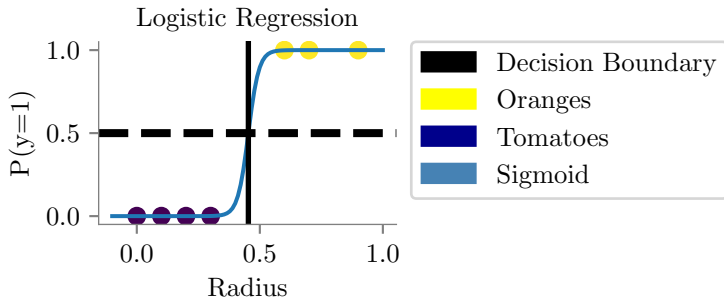
- Have a decision function similar to the above (but not so sharp and discontinuous)

Ideal boundary



- Have a decision function similar to the above (but not so sharp and discontinuous)
- Aim: use linear regression still!

Idea: Use Linear Regression



Question. Can we still use Linear Regression?

Answer. Yes! Transform $\hat{y} \rightarrow [0, 1]$

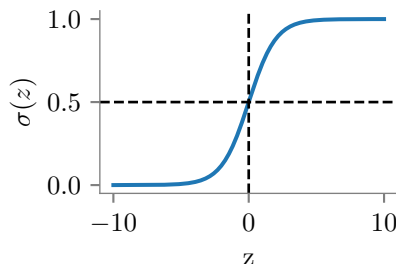
Logistic / Sigmoid Function

$$\hat{y} \in (-\infty, \infty)$$

$\phi = \text{Sigmoid / Logistic Function } (\sigma)$

$$\phi(\hat{y}) \in [0, 1]$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Logistic / Sigmoid Function

$$z \rightarrow \infty$$

Logistic / Sigmoid Function

$$z \rightarrow \infty$$

$$\sigma(z) \rightarrow 1$$

Logistic / Sigmoid Function

$$Z \rightarrow \infty$$

$$\sigma(Z) \rightarrow 1$$

$$Z \rightarrow -\infty$$

Logistic / Sigmoid Function

$$Z \rightarrow \infty$$

$$\sigma(Z) \rightarrow 1$$

$$Z \rightarrow -\infty$$

$$\sigma(Z) \rightarrow 0$$

Logistic / Sigmoid Function

$$Z \rightarrow \infty$$

$$\sigma(Z) \rightarrow 1$$

$$Z \rightarrow -\infty$$

$$\sigma(Z) \rightarrow 0$$

$$Z = 0$$

Logistic / Sigmoid Function

$$Z \rightarrow \infty$$

$$\sigma(Z) \rightarrow 1$$

$$Z \rightarrow -\infty$$

$$\sigma(Z) \rightarrow 0$$

$$Z = 0$$

$$\sigma(Z) = 0.5$$

$$P(y = 0|X) = 1 - P(y = 1|X) = 1 - \frac{1}{1 + e^{-\mathbf{x}\theta}} = \frac{e^{-\mathbf{x}\theta}}{1 + e^{-\mathbf{x}\theta}}$$

$$P(y = 0|X) = 1 - P(y = 1|X) = 1 - \frac{1}{1 + e^{-\mathbf{x}\boldsymbol{\theta}}} = \frac{e^{-\mathbf{x}\boldsymbol{\theta}}}{1 + e^{-\mathbf{x}\boldsymbol{\theta}}}$$

$$\therefore \frac{P(y = 1|X)}{1 - P(y = 1|X)} = e^{\mathbf{x}\boldsymbol{\theta}} \implies \mathbf{x}\boldsymbol{\theta} = \log \frac{P(y = 1|X)}{1 - P(y = 1|X)}$$

Why? Squared loss + sigmoid creates non-convex surface:

- Sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$ is non-linear

Why? Squared loss + sigmoid creates non-convex surface:

- Sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$ is non-linear
- Composition $(\sigma(\mathbf{X}\boldsymbol{\theta}) - y)^2$ has multiple local minima

Why? Squared loss + sigmoid creates non-convex surface:

- Sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$ is non-linear
- Composition $(\sigma(\mathbf{X}\boldsymbol{\theta}) - y)^2$ has multiple local minima

Why? Squared loss + sigmoid creates non-convex surface:

- Sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$ is non-linear
- Composition $(\sigma(\mathbf{X}\boldsymbol{\theta}) - y)^2$ has multiple local minima
- No guarantee gradient descent finds global optimum

Why? Squared loss + sigmoid creates non-convex surface:

- Sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$ is non-linear
- Composition $(\sigma(\mathbf{X}\theta) - y)^2$ has multiple local minima
- No guarantee gradient descent finds global optimum
- **This is why we need cross-entropy loss instead!**

This cost function is called cross-entropy.

This cost function is called cross-entropy.
Why?

What is the interpretation of the cost function?

What is the interpretation of the cost function?

Let us try to write the cost function for a single example:

What is the interpretation of the cost function?

Let us try to write the cost function for a single example:

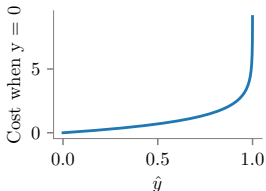
$$J(\theta) = -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)$$

What is the interpretation of the cost function?

Let us try to write the cost function for a single example:

$$J(\theta) = -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)$$

First, assume y_i is 0, then if \hat{y}_i is 0, the loss is 0; but, if \hat{y}_i is 1, the loss tends towards infinity!



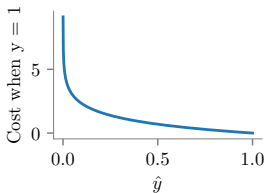
What is the interpretation of the cost function?

$$J(\theta) = -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)$$

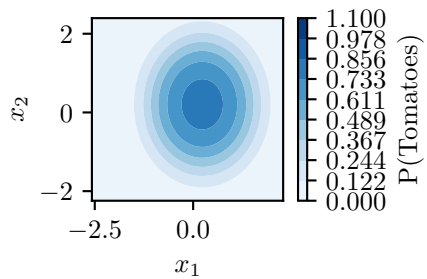
What is the interpretation of the cost function?

$$J(\theta) = -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)$$

Now, assume y_i is 1, then if \hat{y}_i is 0, the loss is huge; but, if \hat{y}_i is 1, the loss is zero!



Bias!



How would you learn a classifier? Or, how would you expect the classifier to learn decision boundaries?

1. Use one-vs.-all on Binary Logistic Regression

1. Use one-vs.-all on Binary Logistic Regression
2. Use one-vs.-one on Binary Logistic Regression

1. Use one-vs.-all on Binary Logistic Regression
2. Use one-vs.-one on Binary Logistic Regression
3. Extend Binary Logistic Regression to Multi-Class Logistic Regression

1. Learn $P(\text{setosa (class 1)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_1)$

1. Learn $P(\text{setosa (class 1)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_1)$
2. $P(\text{versicolor (class 2)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_2)$

1. Learn $P(\text{setosa (class 1)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_1)$
2. $P(\text{versicolor (class 2)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_2)$
3. $P(\text{virginica (class 3)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_3)$

1. Learn $P(\text{setosa (class 1)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_1)$
2. $P(\text{versicolor (class 2)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_2)$
3. $P(\text{virginica (class 3)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_3)$
4. Goal: Learn $\theta_i \forall i \in \{1, 2, 3\}$

1. Learn $P(\text{setosa (class 1)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_1)$
2. $P(\text{versicolor (class 2)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_2)$
3. $P(\text{virginica (class 3)}) = \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_3)$
4. Goal: Learn $\theta_i \forall i \in \{1, 2, 3\}$
5. Question: What could be an \mathcal{F} ?

1. Question: What could be an \mathcal{F} ?

1. Question: What could be an \mathcal{F} ?
2. Property: $\sum_{i=1}^3 \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_i) = 1$

1. Question: What could be an \mathcal{F} ?
2. Property: $\sum_{i=1}^3 \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_i) = 1$
3. Also $\mathcal{F}(\mathbf{z}) \in [0, 1]$

1. Question: What could be an \mathcal{F} ?
2. Property: $\sum_{i=1}^3 \mathcal{F}(\mathbf{X}\boldsymbol{\theta}_i) = 1$
3. Also $\mathcal{F}(\mathbf{z}) \in [0, 1]$
4. Also, $\mathcal{F}(\mathbf{z})$ has squashing properties: $R \mapsto [0, 1]$

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$
 $= -(0 \times \log(0.1) + 1 \times \log(0.8) + 0 \times \log(0.1))$

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$
 $= -(0 \times \log(0.1) + 1 \times \log(0.8) + 0 \times \log(0.1))$
Tends to zero

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$
 $= -(0 \times \log(0.1) + 1 \times \log(0.4) + 0 \times \log(0.1))$

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$
= $-(0 \times \log(0.1) + 1 \times \log(0.4) + 0 \times \log(0.1))$
High number! Huge penalty for misclassification!

More generally,

More generally,

$$J(\theta) = -\left\{ \sum_{i=1}^N y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right\}$$

More generally,

$$J(\theta) = -\left\{ \sum_{i=1}^N y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right\}$$

$$J(\theta) = -\left\{ \sum_{i=1}^N y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right\}$$

Extend to K-class:

$$J(\theta) = -\left\{ \sum_{i=1}^N \sum_{k=1}^K y_i^k \log(\hat{y}_i^k) \right\}$$

What is the key difference between sigmoid and softmax functions?

What is the key difference between sigmoid and softmax functions?

What is the key difference between sigmoid and softmax functions?

Why do we use cross-entropy loss instead of squared error?

What is the key difference between sigmoid and softmax functions?

Why do we use cross-entropy loss instead of squared error?

What is the key difference between sigmoid and softmax functions?

Why do we use cross-entropy loss instead of squared error?

How does regularization help in logistic regression?

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function
- **Linear Decision Boundary:** Creates linear separation in feature space

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function
- **Linear Decision Boundary:** Creates linear separation in feature space

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function
- **Linear Decision Boundary:** Creates linear separation in feature space
- **Maximum Likelihood:** Optimized using gradient-based methods

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function
- **Linear Decision Boundary:** Creates linear separation in feature space
- **Maximum Likelihood:** Optimized using gradient-based methods

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function
- **Linear Decision Boundary:** Creates linear separation in feature space
- **Maximum Likelihood:** Optimized using gradient-based methods
- **Cross-Entropy Loss:** Appropriate for classification problems

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function
- **Linear Decision Boundary:** Creates linear separation in feature space
- **Maximum Likelihood:** Optimized using gradient-based methods
- **Cross-Entropy Loss:** Appropriate for classification problems

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function
- **Linear Decision Boundary:** Creates linear separation in feature space
- **Maximum Likelihood:** Optimized using gradient-based methods
- **Cross-Entropy Loss:** Appropriate for classification problems
- **No Closed Form:** Requires iterative optimization (gradient descent)

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function
- **Linear Decision Boundary:** Creates linear separation in feature space
- **Maximum Likelihood:** Optimized using gradient-based methods
- **Cross-Entropy Loss:** Appropriate for classification problems
- **No Closed Form:** Requires iterative optimization (gradient descent)

Key Takeaways

- **Probabilistic Model:** Outputs probabilities via sigmoid function
- **Linear Decision Boundary:** Creates linear separation in feature space
- **Maximum Likelihood:** Optimized using gradient-based methods
- **Cross-Entropy Loss:** Appropriate for classification problems
- **No Closed Form:** Requires iterative optimization (gradient descent)
- **Regularization:** L1/L2 help prevent overfitting